

Supporting Runtime Reconfiguration on Network Processors

Kevin Lee
Computing, Lancaster University
Lancaster, LA1 4WA
leek@comp.lancs.ac.uk

Geoffrey Coulson
Computing, Lancaster University
Lancaster, LA1 4WA
geoff@comp.lancs.ac.uk

Abstract

Network Processors (NPs) are set to play a key role in the next generation of networking technology. They have the performance of ASIC-based routers whilst offering a high degree of programmability. However, the programmability potential of NPs can only be realised with appropriate software. In this paper we argue that specialised software to support runtime reconfiguration is needed to fully exploit the potential of NPs. We first justify supporting runtime reconfiguration on NPs by offering real-world scenarios and discussing the issues associated with these. We then demonstrate how runtime reconfiguration can be achieved in practice through a case study of our component-based programming approach on the Intel IXP2400 NP.

1. Introduction

Network Processors (NPs) are multiprocessor-based hardware units that have the ability to perform relatively complex packet processing tasks in software at line speeds when compared to contemporary devices. They typically consist of a set of heterogeneous processors including packet processors, dedicated devices such as encryption engines, general purpose processors, and a high-speed interconnect [1].

NPs can be seen as an attempt by hardware vendors to fulfill the growing need for network hardware elements that support high throughput while also offering increased programmability. Programmability is seen as crucial in supporting system evolution so that new protocols and services can be accommodated without designing new hardware. Their programmability makes NPs very widely applicable—e.g. they are being applied in networked devices, as edge-network routers and even in the network core [2].

In addition, it is now becoming recognised [3] that *runtime reconfiguration* is a desirable characteristic of software for NPs. Runtime reconfiguration is useful for a number of applications, including dynamically extensible services [4],

network resource management [5], configurable network-based encryption [6], and offloading of processing [7]. In addition the active networking (AN) community have been heavily involved in investigating the use of NPs in their field [8]. This is because ANs require significant data-plane processing and also require routers to expose their state to allow reconfiguration of forwarding functions.

The aim of the research discussed in this paper is to investigate the potential and benefits of runtime reconfiguration in NPs. Our research focuses on the provision of generic mechanisms that can potentially be applied in a wide range of scenarios including all of the above. We adopt a runtime component-based approach in which fine-grained components on the NP can be dynamically (un)loaded and (dis)connected in a principled manner. In this paper we illustrate the generality of our approach by focusing on applying it in a set of representative scenarios. We use the Intel IXP 2400 [9] as representative of the state of the art of the current generation of NPs. We also argue that a flexible runtime reconfiguration capability need not be bought at the expense of performance.

The remainder of this paper is structured as follows. In section 2 we motivate runtime reconfiguration by outlining common reconfiguration scenarios and then provide background on the Intel IXP2400. Section 3 then presents details of OpenCOM, our runtime reconfiguration capable programming platform, and its deployment on the IXP2400. In section 4 we show how the reconfiguration scenarios can be realised by OpenCOM. Finally, we discuss related work in section 5 and in section 6 offer our conclusions.

2. Background

2.1. Runtime Reconfiguration

2.1.1. Dynamic Proxying. In general terms, proxying is a technique for allowing clients and devices to make indirect connections to network services via a shared intermediary. It is used both to limit the network load incurred in pro-

viding access to external network resources and to provide value-added services. The proxy notion can be applied in a wide range of settings including web caching, VoIP proxying, and media transcoding. Furthermore, it can involve a range of generic techniques including combining client requests, diverting connections, denying connections, or creating encrypted tunnels.

Currently, proxying is typically performed at the network edge on dedicated devices. However, with NPs it becomes possible to deploy proxies on routers inside the network. Furthermore, such proxies could be deployed on the fly and on demand. To support such dynamic proxies, a NP would need a software framework that incorporated an extensible classifier to identify specific flows, plus the ability to instantiate proxy components depending on the service required. The benefits would be a minimisation of latency as well as a maximisation of flexibility. Deployment overhead could also be minimised by caching proxies on the NP.

2.1.2. Adaptive Load Balancing. On standard network routers, flows are either not differentiated or are differentiated in a relatively static manner (e.g. using diffserv). There is no capability to adapt the resources dedicated to different flows in a fine grained manner depending on current application needs or traffic patterns.

With NPs, however, it is possible to dynamically deploy resources to different flows. For example, a given number of hardware threads, or packet processors, could be dedicated to high-priority VoIP flows, depending on patterns of demand (e.g. as a function of the time of day). As well as packet forwarding, this also applies to per-flow processing such as in-band transcoding. Furthermore, because NPs have the ability to process and forward traffic while simultaneously analysing the traffic to determine a suitable load balancing policy, they offer the ability to perform “intelligent” load balancing. In contrast to off-line load balancing, fine-grained load-balancing mechanisms can range from diverting flows to different routes to replicating processing/classification code across multiple packet processors.

2.2. The Intel IXP2400

The Intel IXP2400 NP [9] consists of a single embedded RISC processor (an Intel XScale), and eight packet processors called “microengines”. It provides a fast bus for communication between its microengines, MAC ports and memory. It also provides shared registers and a range of memory types (i.e. SRAM, SDRAM). In addition, it provides ‘next-neighbour’ registers that provide a dedicated interconnect between two ‘adjacent’ microengines.

The microengines themselves are 233-600MHz CPUs whose instruction set provides for I/O to/from MAC-ports, packet queuing support, and checksumming. They support

hardware threads with zero context switch overhead and can be programmed either in assembler or C.

In normal operation, the IXP2400 uses the microengines to support the data plane and the more general XScale to support the control plane. The shared registers and memory are typically used together at the software level to realise inter-processor communication.

3. OpenCOM

3.1. Programming Model

OpenCOM [10] is a language independent component-based programming platform for building low-level systems software. The core principles of OpenCOM are *components*, *capsules*, *caplets*, *interfaces*, *receptacles*, and *bindings*.

- **Components, capsules and caplets** *Components* are encapsulated units of functionality and deployment that interact with their environment (i.e. other components) exclusively through interfaces and receptacles. Components carry negligible inherent overhead and can effectively be used in extremely fine grained compositions. Crucially, OpenCOM is a *runtime* component model meaning that (unlike, say, NP-Click [11]) components can be dynamically deployed at any time during run-time. The locus of component deployment is either a *capsule* or a *caplet*; the latter are subscopes of the former. Different caplets can also host components written in different *component styles*. Component styles are different system-level implementations of components which may have different representations and different semantics (e.g. because they run on different CPU types). Accommodating heterogeneous component styles enables OpenCOM to transparently support multiple deployment environments in the same capsule environment.

Each capsule offers a simple run-time API for component life-cycle management (i.e. loading components into the capsule and instantiating and destroying them), and interface/receptacle binding (see below). To accomplish loading, the model supports the notion of *plug-in loaders*. New loaders with different behaviours can be added at runtime, and they can be selected according to their particular properties. The loading of components into a capsule can be requested by any component hosted by the capsule no matter which caplet it resides in (this is referred to as *third-party deployment*).

- **Interfaces and receptacles** *Interfaces* are units of service provision offered by components; they are expressed in terms of sets of operation signatures and associated datatypes. For language independence, OMG IDL is used as a specification language (note that this does *not* imply an overhead of CORBA-like stubs and skeletons!). Components can support multiple interfaces: this is useful in recognising separations of concerns (e.g. between base functionality and management). *Receptacles* are ‘anti-interfaces’ used

to make explicit the dependencies of components on other components. Receptacles are key to supporting a third-party style of composition (to complement the third-party deployment referred to above): when third-party-deploying a component into a capsule, one knows by looking at the component's receptacles precisely which other component types must be present to satisfy its dependencies.

- **Bindings** Finally, *bindings* are associations between a single interface and a single receptacle that reside in a common capsule (but not necessarily a common caplet). Similarly to plug-in loaders, OpenCOM also supports a notion of *plug-in binders*. The idea is to give access to a range of binding mechanisms with varying characteristics. As mentioned, the creation of bindings is inherently third-party in nature; it can be performed by any party within the capsule (i.e. not only by the 'first-party' components whose interface or receptacle participates in the binding).

3.2. Higher-level abstractions

Above the granularity of individual components, a key pattern employed in OpenCOM programming is to construct applications or systems in terms of *component frameworks* (CFs). CFs are tightly-coupled sets of components that work together to address some specific area of functionality. They accept 'plug-in' components, deployed at runtime, which somehow modify the CF's behaviour. CFs also impose constraints on their plug-ins to guard against nonsensical compositions. As an example, consider a "protocol stack" CF which accepts protocol components as plug-ins, and constrains these plug-ins to be composed linearly. A CF that we employ specifically in NP environments, the Router CF, is discussed in section 3.4.

We also support a number of generic services that facilitate the construction of complex systems. These are themselves implemented in terms of components and are thus optional in any given capsule configuration. Key among these is a set of *reflective meta-models* [10] that facilitate dynamic reconfiguration of systems by permitting different system aspects to be inspected, adapted and extended at runtime. As examples: the *architecture meta-model* exposes the compositional topology of the components in a capsule in terms of a causally-connected graph structure; the *interface meta-model* allows one to discover information about interface types at runtime and to invoke interface instances that are dynamically discovered at runtime; and the *interception meta-model* allows one to interpose interceptors at bindings between component interfaces.

3.3. OpenCOM on the Intel IXP2400

We now consider how the above-described OpenCOM concepts are applied in NPs such as the IXP2400. First, the scoping-related notions of capsules and caplets are useful

in distinguishing different processors and types of processors on the NP in a generic manner. Thus we map a single capsule onto the entire NP, and sub-scope individual microengines, and the XScale control processor, as caplets. The OpenCOM runtime runs in the XScale caplet; all the other caplets are 'slaves' of this 'central' runtime and incur only minimal memory overheads. The memory footprint of the central runtime itself is of the order of 300KB.

Microengine caplets are implemented on the bare microengine hardware. The notion of caplets is also useful in isolating untrusted code, which is important in active networking environments. For example, a Java sandbox could be isolated as a caplet on the XScale.

The pluggable loader concept is closely associated with capsules/caplets. Typically, at least one loader is provided for each type of caplet, and each will know how to load components into the hardware environment underlying its particular caplet type. For example, we employ one loader for the XScale caplet and another for the microengine caplets. Importantly, the OpenCOM API allows selective transparency in the use of loaders. If full loader-selection transparency is desired, one can issue a call of the form *load(component_c1, caplet_1)* which will deduce an appropriate loader type from meta-data attached to *component_c1*, and use this to load the component into the designated caplet. This masks the fact that different components may be implemented in different machine languages. Alternatively, one can opt for complete control and zero transparency by issuing a call of the form *load(component_c1, caplet_1, loader_3)*.

The pluggable binder concept is equally central to the component model's abstraction power. If we don't know or care in which caplets our two target components are located, we can say *bind(interface_1, receptacle_15)* and an appropriate loader will be selected according to location-related meta-data attached to the components that own the specified interface and receptacle. On the other hand, if it is important to select a particular mechanism, we can say *bind(interface_1, receptacle_15, loader_4)*. And so on.

A final crucial property of the component model is its radically third-party nature in terms of loading and binding. Thanks to this, a component on a microengine can load and bind two components on the XScale, and a component on the XScale can load and bind microengine components using exactly the same syntax as if it were dealing with local XScale components.

As an example of this, and of OpenCOM's ability to abstract over heterogeneity of the IXP2400, consider the following pseudo-code segment:

```
template mtemp, xtemp;
comp_inst mcaplet, mloader, xcaplet, xloader,
        mcomp, xcomp, binding, cbinder;
ipnt_inst iface, recpt;

// load and instantiate the components
xtemp = load(XSCALECOMP1, xloader, xcaplet);
```

```

mtemp = load(MICROCOMP1, mloader, mcaplet);
xcomp = instantiate(xtemplate);
mcomp = instantiate(mtemplate);

// bind the components using a cross-caplet binder
binding = bind(xcomp.iface, mcomp.recpt, cbinder);

```

This example assumes that two caplets have been established: *xcaplet* is a caplet on the XScale and *mcaplet* is a caplet on one of the microengines. The code loads and instantiates two components, one in each caplet, and then binds the two using a cross-caplet binder. Of course, programming would normally be done at the level of component frameworks which raises the level of abstraction still further; but this simple example shows the abstraction power of OpenCOM even at its lowest level.

As examples, we now briefly describe example loader and binder plug-ins that are associated with the microengine caplets:

- The *microengine loader plug-in* provides the illusion of dynamic loading despite the fact that the microengine hardware only allows modification of its instruction store when the CPU is stopped [12]. The basic capability provided by the microengine hardware is to i) stop the microengine, ii) read/ write arbitrary instruction store locations, and then iii) restart it at a hard-wired address. To achieve transparent dynamic loading it is therefore necessary for the loader to not only load the new component but also to patch the (hard-wired) restart address so that subsequent execution resumes at the point it left off. The loader also has the ability to autonomously move code around within the instruction store to avoid memory fragmentation as components are loaded and unloaded.

- Our *intra-microengine binder plug-in* is strongly coupled to the loader just described. It builds on the NetBind-pioneered technique of creating bindings by ‘morphing’ jump instructions [13]. Together with the loader discussed above, the binder supports multiple instantiations of components (NetBind only supports singleton components). The single argument and return value are passed via a designated register. The necessary entry and exit point information is obtained from IDL meta-data attached to the packaged component, which is transformed from relative offsets to absolute offsets by the loader. The overhead of a binding created by this binder in calling a null operation with no arguments or return values is only five (1-cycle) instructions. These subsume passing on the stack a pointer to the per-instance state vector of the called component, and the return address.

3.4. Router CF

We have designed a “Router CF” which accepts, as plug-ins, OpenCOM components that perform arbitrary, user-defined packet-forwarding functions (we also provide “standard” components that interface to network cards and

wrap efficient kernel-user space communications mechanisms). All components are required to conform to the following rules, which are checked by the CF at run-time when the component is loaded:

- compliant components must support appropriate numbers and combinations of specific packet passing interfaces/receptacles (called *IPacketPush* and *IPacketPull*): these respectively enable push- and pull- oriented inter-component communication); it is possible to dynamically add/ remove instances of these interfaces as long as the CF’s rules remain satisfied
- compliant components may (optionally) support an *IClassifier* interface which exports an operation *register_filter()* that is used to install packet-filters; if *IClassifier* is supported, the component must honour the semantics of installed filter specifications in terms of the particular named outgoing *IPacketPush* and *IPacketPull* interface(s) on which each incoming packet should be emitted;
- compliant components may be *composite*, in which case all their internal constituents must (recursively) conform to the CF’s rules; additionally composite components should contain a so-called *controller* component that manages the other internal constituents.

4. Runtime Reconfiguration on the Intel IXP2400

In this section we illustrate how the runtime reconfiguration scenarios illustrated in 4.1 can be achieved using OpenCOM on the IXP2400. We then discuss the performance implications of our implementation in section 4.2.

4.1. Realising the Scenarios

4.1.1. Dynamic Proxying. We have implemented a dynamic in-band transcoding scenario that is straightforwardly built on top of the programmable classifier discussed above. In-band transcoding (e.g. MPEG) has not been possible in the past without intercepting the multimedia traffic and off-loading it to a separate server. As noted earlier, this requires introducing new hardware into the system and also potentially increases the end-to-end latency of media transmission.

In our implementation, a dynamically-deployed manager component, residing in the XScale caplet, programs the classifier to detect flows of interest (i.e. as designated by out-of-band control messages from an application). On learning of these, the manager component loads and instantiates a suitable transcoder on the XScale (or a microengine as appropriate). On doing so, the router CF selects an appropriate loader and checks the loaded component (using the interface meta-model) for conformance to its rules. The manager then instructs the classifier to bind to this component and to forward the flow to it. Finally, it uses the architecture meta-

model to locate the forwarder and bind the transcoder to it using an appropriate binder.

This implementation is illustrated in Figure 1 which shows the state of a router with the transcoder manager deployed. It shows the reconfigured area (within the dashed line) containing the manager on the XScale and transcoders on both the XScale and the microengines. Note that the microengines can only support primitive transcoders such as frame-droppers; note further, though, that the programming model makes it as straightforward to deploy a transcoder on a microengine as on the XScale.

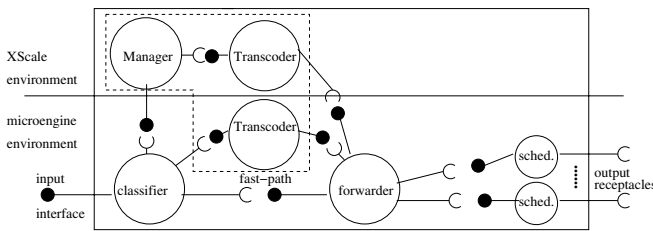


Figure 1. Transcoding on the Intel IXP2400

4.1.2. Adaptive Load Balancing. The options for load balancing on the Intel IXP2400 are numerous. In typical operation, the bulk of packets traversing the IXP2400 are processed and forwarded by the microengines. At different times in the lifecycle of a typical deployment the load on particular microengines will be increased or decreased. Therefore, to balance increased packet load in the IXP2400 one of the options we have is to replicate packet processing code on additional microengines.

Figure 2 shows the placement of components after a simple method of load balancing has been performed. Before the load-balancing is performed, the top four components constitute the deployment in the microengine. Load-balancing is performed by replicating the “IPv4 header processing” and “forwarding” components on additional microengines. As can be seen from the diagram, the classifier is load balancing across the two chains of components. The dashed line indicates the reconfigured area. This style of load-balancing would be appropriate when there is a significant increase of a type of packet flow which needs additional processing by the NP.

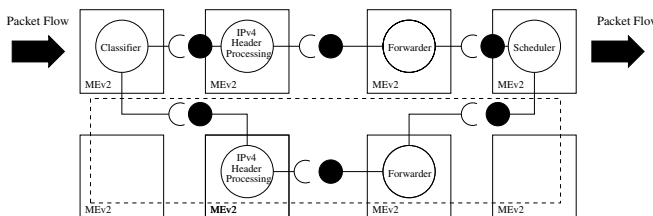


Figure 2. Load Balancing using Packet Processors

This scenario shows that, based on the network situation, the NP control processor can add and remove components to alter the processing capacity of the NP. NPs generically contain a number of packet processors, which perform the majority of in-band packet processing [14]. Therefore this style of load balancing should be applicable in other NP architectures.

4.2. Performance

One of the main determining factors in the acceptance of support for runtime reconfiguration on NPs is the overhead incurred. This breaks down into two aspects: i) the overhead of actually performing reconfiguration operations; and ii) the inherent overhead of potentially-reconfigurable software. We discuss these in turn.

The major determinant of the overhead of reconfiguration operations on the IXP2400 is that the microengines need to be stopped before new code can be loaded. Our measurements show that to halt, update and restart a microengine takes a total time of 60ms. Our IXP2400 development board contains 3 OC-48 ports which can each deliver 2.488Gbps, a total of 7.464Gbps. A delay in the system would therefore require a maximum theoretical of 56MB to buffer all incoming packets and avoid dropping packets, well within the means of a NP. Furthermore, wholesale reconfigurations of all the microengines would be relatively uncommon.

The second factor to be considered is the inherent overhead of potentially-reconfigurable component-based software; mainly attributed to the bindings between components. To evaluate the throughput overhead of instantiating OpenCOM components on the IXP2400, we utilised two 3COM 3C996-SX network interface cards which were used to send and receive packets through a Radysis ENP-2611 which consists of a IXP2400 NP and 3X 2.5Gbps fibre ports. The XScale CPU of the IXP2400 was bootstrapped with Linux 2.6.11 and all microengine code was loaded and bound from an OpenCOM instantiation on the XScale CPU. The following throughput results were collected using the Thulay tool which uses a client/server approach to measure TCP throughput.

A single OpenCOM component which performed a simple layer 2 bridging operation between two fiber channel connections was deployed on a single thread of a microengine. The component was capable of processing packets at a sustained rate of 632.42Mbps end-to-end compared to a monolithic Intel implementation at 632.81Mbps. Additional ‘null’ components were then instantiated directly in the data-path between the send and receive portions of the bridging component.

Table 1 presents end-to-end throughput and latency figures for different numbers of ‘null’ OpenCOM components instantiated on the IXP2400 as described. It shows that the overhead of inserting five or less OpenCOM microengine

Table 1. Throughput and Latency of Components

Number of Components	Throughput	Latency
Intel Implementation	632.81Mbps	0.52ms
1 Component	632.42Mbps	0.54ms
5 Components	614.86Mbps	0.59ms
10 Components	603.25Mbps	0.64ms
15 Components	589.82Mbps	0.88ms
20 Components	567.39Mbps	1.17ms
50 Components	496.03Mbps	2.81ms
100 Components	435.72Mbps	3.65ms

components is minimal, inserting between 10 and 20 components introduces a sizable lag into the system, this might be considered acceptable for the advantages offered. The insertion of ten or more components introduces a sizable lag into the system which would be considered unacceptable for a high-speed router. In addition the figures for latency correlate with the throughput figures with increasing latency with more components. The implication of these results is that the most effective way to deploy OpenCOM components is using multiple short pipelines of five or less components.

5. Related work

Intel's *MicroACE* [12] is an NP-based programming platform targeted at the IXP2400 and other Intel IXA products. The *MicroACE* model is that proxy-like software elements (called *active computing elements* or ACEs) on the IXP2400's XScale control processor are 'mirrored' by blocks of code (called microblocks) that run on micro-engines. Although it provides a useful degree of abstraction, the *MicroACE* model is *static* in nature. It does not support any type of runtime reconfiguration.

NetBind [13] provides the abstraction of a set of packet-processing components that can be bound into a data path. This is done by adopting the convention of a standard entry and exit instruction sequence for microblocks, and offering the capability to dynamically 'morph' jump instructions in these sequences so that execution is transferred to the entry point of the microblock to be executed next. However, it offers no abstraction over the NP's memory organisation, interconnects or processors and therefore offers no more design portability across different NPs than *MicroACE*.

NP-Click [11] is another component-based programming model targeted at the IXP1200. Its components have typed *ports*; and connections between ports can be designated as either 'push' or 'pull' which provides declarative control over flow of control and threading. *NP-Click* falls short of providing a generic approach to NP programming. While it abstracts particular features of the IXP1200, it has no notion of abstracting arbitrary architectures in a principled way, and thereby encouraging design portability and transferable skills across NP types. In addition, *NP-Click* provides no support for dynamic reconfiguration.

6. Conclusions

We have argued that developing NP software with support for runtime reconfiguration enables the full potential of NPs to be realised, and that this yields significant benefits for high-speed routing platforms. More specifically, we have introduced a number of runtime reconfiguration scenarios for NP platforms and showed how they can be implemented on the Intel IXP2400 using our OpenCOM programming platform.

We also argue the approach we outline is in principle applicable not only to the IXP2400, but to a range of NP architectures. This claim is made on the basis of the generality of the OpenCOM platform as discussed in 3.3 and on the basis of a study of the mapping of OpenCOM to other NP architectures [14].

References

- [1] D. Comer. Network Systems Design using Network Processors, IXP edition. 2003.
- [2] Heavy Reading. Network processors: A heavy reading competitive analysis. In *Vol. 3, No. 2*, January 2005.
- [3] I.A. Troxel, A.D. George, and S. Oral. Design and analysis of a dynamically reconfigurable network processor. In *IEEE Conference on Local Computer Networks*, Nov 2002.
- [4] L. Ruf, R. Keller, and B. Plattner. A Scalable High-performance Router Platform Supporting Dynamic Service Extensibility On Network and Host Processors. In *IEEE Conference on Pervasive Services*, Jul 2004.
- [5] A. Gavrilovska, K. Schwan, and O. Nordstrom. Network processors as building blocks in overlay networks. In *11th Symposium on High Performance Interconnects*, Aug 2003.
- [6] S. Harper. Phd thesis: A Secure Adaptive Network Processor, Bradley Department of Electrical and Computer Engineering Blacksburg, Virginia, April 2003.
- [7] C. Lee et al. Software/hardware reconfigurable network processor for space networks. In *MAPLS 2001*, Jan 2001.
- [8] A. Kind, R.Pletka, and M.Waldvogel. The role of network processors in active networks. In *IWAN 2003*, Japan, 2003.
- [9] Intel Corporation. Intel IXP2400 Network Processor. In *Datasheet 301164-011*. Intel Corporation, Feb 2004.
- [10] G. Coulson, G. Blair, P. Grace, A. Joolia, K. Lee, and J. Ueyama. A component model for building systems software. In *IASTED 2004*, Cambridge, MA, USA, 2004.
- [11] N. Shah, W. Plishker, and K. Keutzer. NP-Click: A Programming Model for the Intel IXP1200. In *2nd Workshop on Network Processors at HPCA-9*, Anaheim, February 2003.
- [12] Intel Corporation. *MicroACE*, Design Document, revision 1.0. Intel Press, Intel Corporation, 2001.
- [13] A. Campbell, M. Kounavis, and D. Villela. NetBind: A Binding Tool for Constructing Data Paths in Network Processor-based Routers. In *IEEE International Conference on Open Architectures*, June 2002.
- [14] K. Lee, G.Coulson, and G. Blair et al. Towards a generic programming model for network processors. In *IEEE International Conference on Networks*, Singapore, Nov 2004.